



CSC2547 Presentation:
Curiosity-driven
exploration

Count-based VS Info gain-based



Sheng Jia, Tinglin Duan
(First year master students)

Planning to Be Surprised: Optimal Bayesian Exploration in Dynamic Environments

Yi Sun, Faustino Gomez, and Jürgen Schmidhuber

VIME: Variational Information Maximizing Exploration

Rein Houthoofd^{§††}, Xi Chen^{††}, Yan Duan^{††}, John Schulman^{††}, Filip De Turck[§], Pieter Abbeel^{††}

Unifying Count-Based Exploration and Intrinsic Motivation

Marc G. Bellemare
bellemare@google.com

Sriram Srinivasan
srsrinivasan@google.com

Georg Ostrovski
ostrovski@google.com

Tom Schaul
schaul@google.com

David Saxton
saxton@google.com

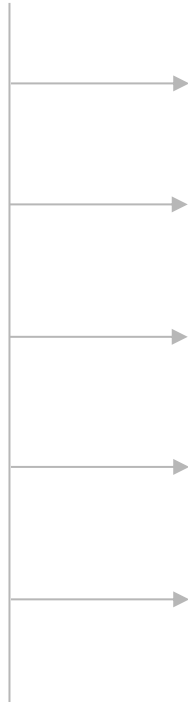
Rémi Munos
munos@google.com

1. PLAN (2011)

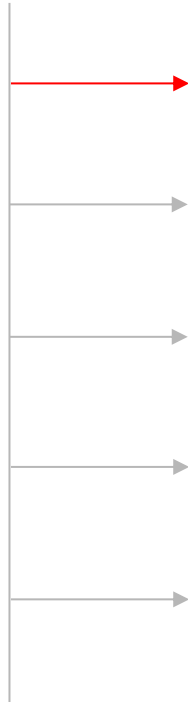
1. VIME (NeurIPS2016)

1. CTS (NeurIPS2016)

Outline

- 
- Motivation, Related Works and Demo
 - Planning to Be Surprised
 - Variational Information Maximizing Exploration
 - Unifying Count-Based Exploration and Intrinsic Motivation
 - Comparisons and Discussion

Outline

- 
- Motivation, Related Works and Demo
 - Planning to Be Surprised
 - Variational Information Maximizing Exploration
 - Unifying Count-Based Exploration and Intrinsic Motivation
 - Comparisons and Discussion

Background

RL+Curiosity

History:

$$\xi_t = (s_1, a_1, s_2, a_2, \dots, s_t)$$

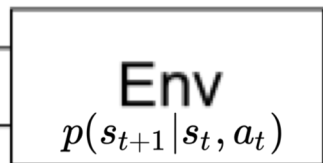
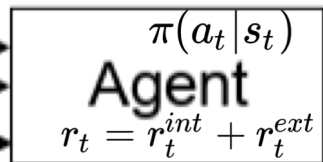
Intrinsic
reward
/exploration
bonus

$$r_t^{int}$$

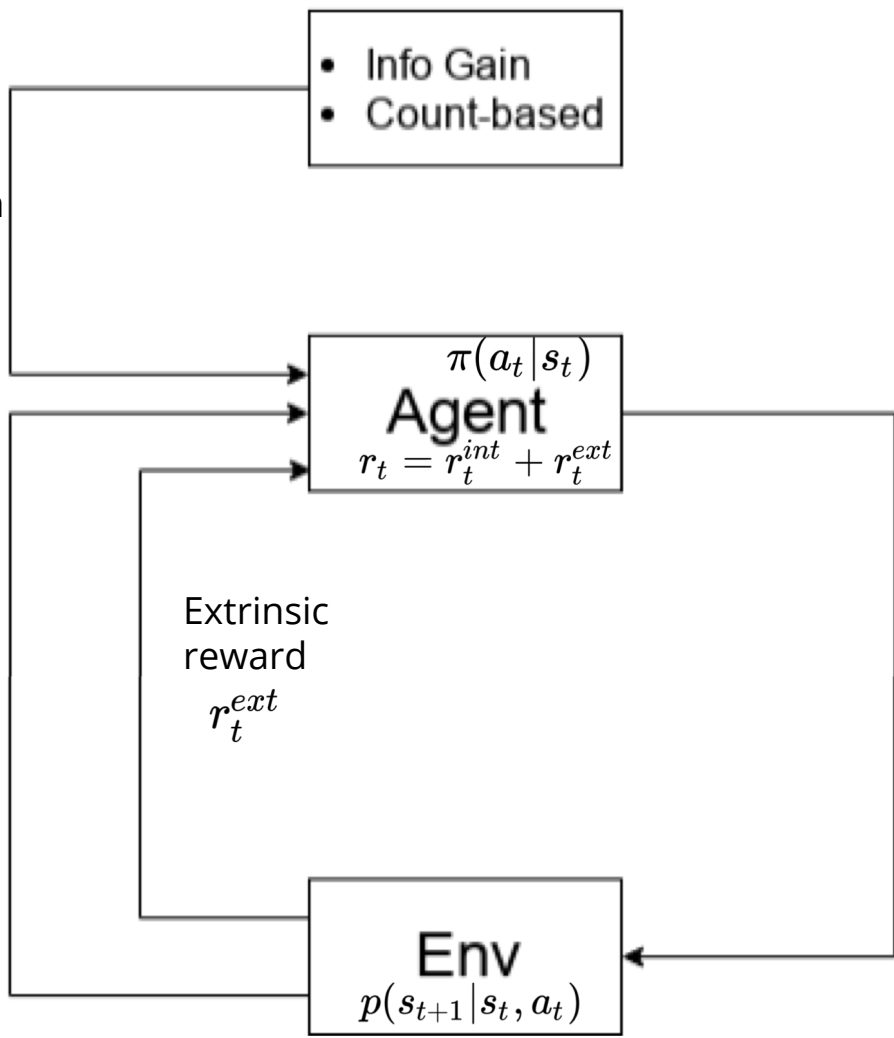
Next
state
 s_{t+1}

Extrinsic
reward
 r_t^{ext}

- Info Gain
- Count-based



action
 a_t



What is exploration?

Intrinsic motivation:

- Reducing the agent's uncertainty over the environment's dynamics.

[Plan] $p(s_{t+1}|\xi_t, a_t; \theta)$ $p(\theta)$

[VIME] $p(s_{t+1}|s_t, a_t; \theta)$

[CTS]

Count-based

- Use (pseudo) visitation counts to guide agents to unvisited states.

Why exploration useful?

DEMO

Sparse Reward Problem
Montezuma's revenge

Our original plot & demo

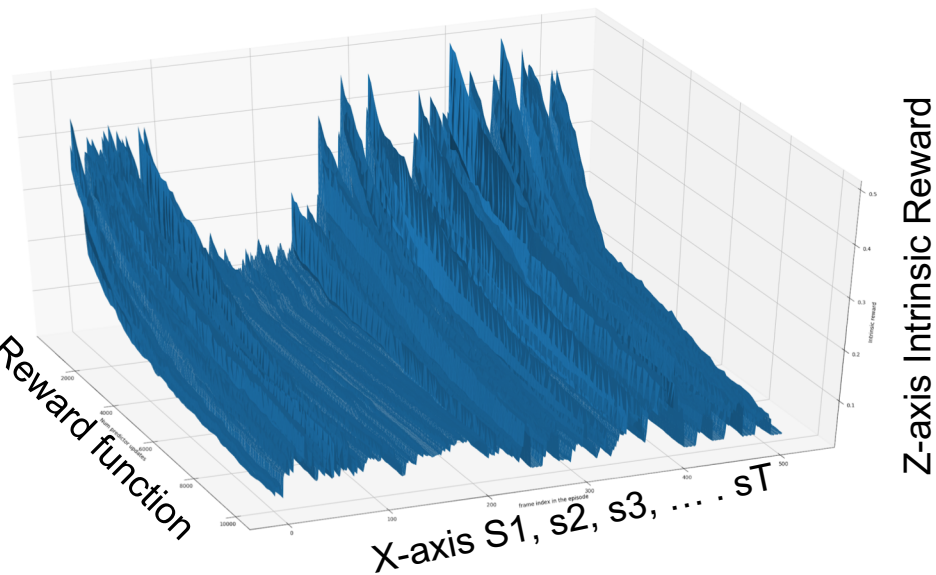
DQN



DQN + Exploration bonus



Y-axis Intrinsic Reward
timestamp
Training Timestamp



Related work (Timeline)

The notion of Intrinsic Motivation

2010 Formal Theory of Creativity, Fun, and Intrinsic Motivation (1990-2010)

L2 prediction error using neural networks

2015 Incentivizing Exploration In Reinforcement Learning With Deep Predictive Models

Pseudocount + Pixel CNN

2017 Count-Based Exploration with Neural Density Models

Pseudocount in 2016 still achieves SOTA for Montezuma's revenge"

2019 On Bonus Based Exploration Methods In The Arcade Learning Environment

2011
PLAN

**Bayesian
Optimal
Exploration**

2016
VIME

**Approximate
"PLAN"**

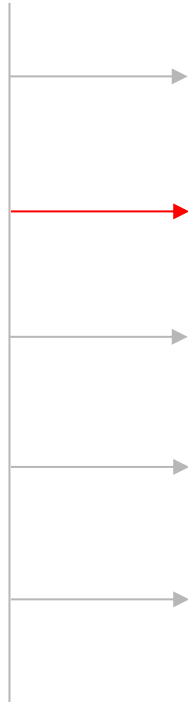
CTS

**Pseudocount
exploration**

2018 Exploration by
Random Network
Distillation

**Distillation error as a
quantification of uncertainty**

Outline

- 
- Motivation, Related Works and Demo
 - Planning to Be Surprised**
 - Variational Information Maximizing Exploration
 - Unifying Count-Based Exploration and Intrinsic Motivation
 - Comparisons and Discussion

[PLAN] contribution

Dynamics model

$$p(s_{t+1} | \xi_t, a_t; \theta)$$

Bayes update for
posterior distribution of
the dynamics model

$$p\left(\underbrace{\theta}_y | \xi_t, a_t, \underbrace{s_{t+1}}_x\right) = \frac{p\left(\underbrace{\theta}_y | \xi_t, a_t\right) p\left(\underbrace{s_{t+1}}_x | \xi_t, a_t; \underbrace{\theta}_y\right)}{p\left(\underbrace{s_{t+1}}_x | \xi_t, a_t\right)}$$

Optimal Bayesian
Exploration based on:

$$q^\tau(\xi_t, a_t) = \mathbb{E}_{s_{t+1} | \xi_t, a_t} [D_{KL} [p(\theta | \xi_t, a_t, s_{t+1}) || p(\theta | \xi_t)]] + \mathbb{E}_{s_{t+1} | \xi_t, a_t} \left[\max_{a_{t+1}} q_\pi^{\tau-1}((\xi_t, a_t, s_{t+1}), a_{t+1}) \right]$$

Expected cumulative
info gain fo tau steps
if performing this
action

Expected one-step info gain

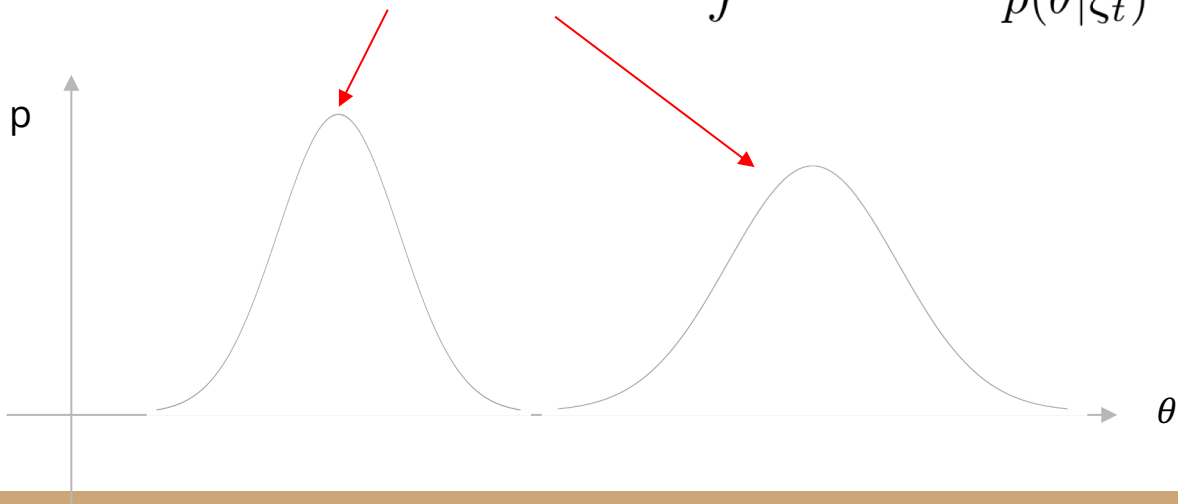
Expected cumulative info
gain for tau-1 steps if
performing this next action

[PLAN] Quantify “surprise” with info gain

$$\xi_t = (s_1, a_1, s_2, a_2, \dots, s_t)$$

$$\xi'_t = (s_1, a_1, s_2, a_2, \dots, s_t, a_t \dots, s_{t'})$$

$$KL(p(\theta|\xi'_t) || p(\theta|\xi_t)) = \int p(\theta|\xi'_t) \log \frac{p(\theta|\xi'_t)}{p(\theta|\xi_t)} d\theta$$



[PLAN] 1-step expected information gain

“1-step expected info gain” “expected immediate info gain”

$$\begin{aligned} & \mathbb{E}_{s_{t+1} \sim P(\cdot | \xi_t, a_t)} [D_{KL} [p(\theta | \xi_t, a_t, s_{t+1}) || p(\theta | \xi_t)]] \\ &= \sum_{s_{t+1}} p(s_{t+1} | \xi_t, a_t) \int p(\theta | \xi_t, a_t, s_{t+1}) \log \frac{p(\theta | \xi_t, a_t, s_{t+1})}{p(\theta | \xi_t)} d\theta \\ &= \sum_{s_{t+1}} \int p(s_{t+1}, \theta | \xi_t, a_t) \log \frac{p(s_{t+1}, \theta | \xi_t, a_t)}{\underbrace{p(\theta | \xi_t)}_{p(\theta | \xi_t, a_t)} p(s_{t+1} | \xi_t, a_t)} d\theta \\ &= I(S_{t+1}; \Theta | \xi_t, a_t) \end{aligned}$$

NOTE: VIME uses this as the Intrinsic reward!

“Mutual info between next state distribution & model parameter”

[PLAN] “Planning to be surprised”

Curious Q-value

Cumulative τ steps info gain

$$q_{\pi}^{\tau}(\xi_t, a_t) = \mathbb{E}_{s_{t+1}, a_{t+1}, s_{t+2}, \dots, a_{t+\tau}, s_{t+\tau+1} | \xi_t, a_t} D_{KL} [p(\theta | \xi_t, a_t, s_{t+1}, a_{t+1}, \dots, a_{t+\tau}, s_{t+\tau+1}) || p(\theta | \xi_t)]$$

Follow a policy

“Planning tau steps”
because not actually
observed yet

Perform an action

[PLAN] Optimal Bayesian Exploration policy

[Method1] Computing optimal curiosity-Q backwards for tau steps

$$q^\tau(\xi_t, a_t) = \mathbb{E}_{s_{t+1}|\xi_t, a_t} [D_{KL} [p(\theta|\xi_t, a_t, s_{t+1}) || p(\theta|\xi_t)]] + \mathbb{E}_{s_{t+1}|\xi_t, a_t} \left[\max_{a_{t+1}} q^{\tau-1}((\xi_t, a_t, s_{t+1}), a_{t+1}) \right]$$

[Method2] Policy Iteration

Repeat applying

Policy evaluation

$$v_\pi^\tau(\xi_t) = \mathbb{E}_{a_t|\xi_t} \left[\mathbb{E}_{s_{t+1}|\xi_t, a_t} [D_{KL} [p(\theta|\xi_t, a_t, s_{t+1}) || p(\theta|\xi_t)]] + v_\pi^{\tau-1}(\xi_t, a_t, s_{t+1}) \right]$$

Policy improvement

$$\pi^\tau(\xi_t) = \arg \max q^\tau(\xi_t, a_t) :$$

[Plan] Non-triviality of curious Q-value

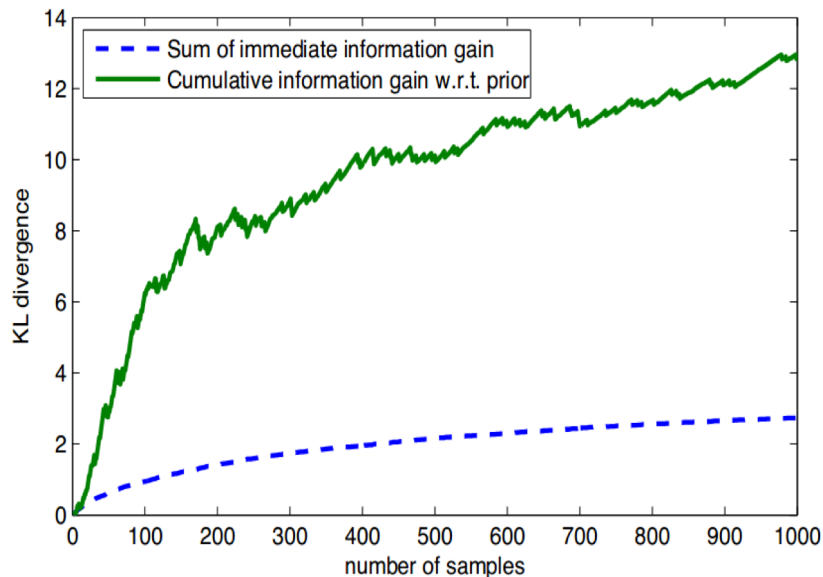
Info gain additive in expectation!

$$\mathbb{E}_{\xi_t'' || \xi_t} [D_{KL} [p(\theta | \xi_t'') || p(\theta | \xi_t)]] = D_{KL} [p(\theta | \xi_t') || p(\theta | \xi_t)] + \mathbb{E}_{\xi_t' || \xi_t} [D_{KL} [p(\theta | \xi_t'') || p(\theta | \xi_t')]]$$

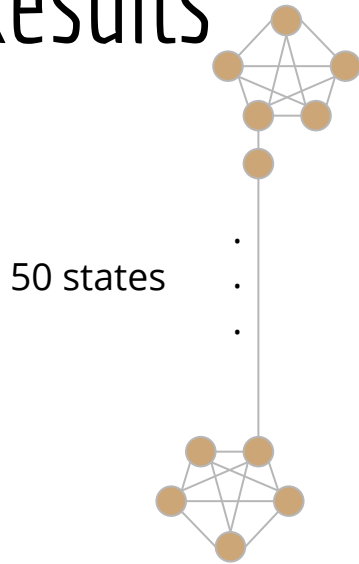
Cumulative \neq Sum

$$D_{KL} [p(\theta | \xi_t'') || p(\theta | \xi_t)] \neq D_{KL} [p(\theta | \xi_t') || p(\theta | \xi_t)] + D_{KL} [p(\theta | \xi_t'') || p(\theta | \xi_t')]$$

Cumulative information gain fluctuates!



[Plan] Results

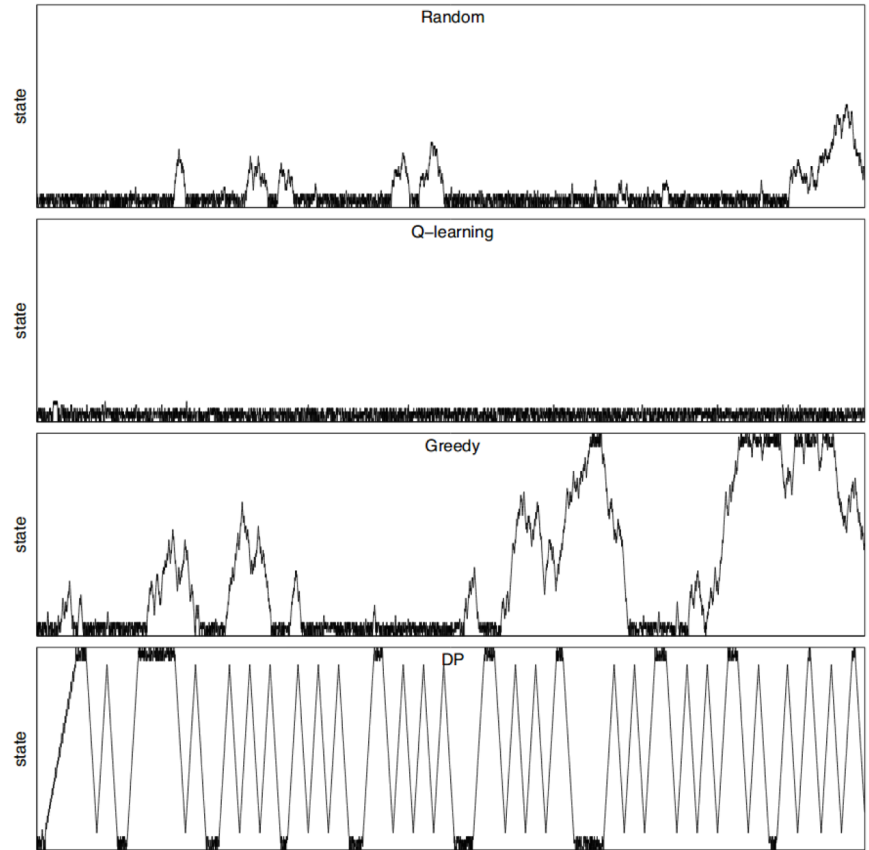


Random

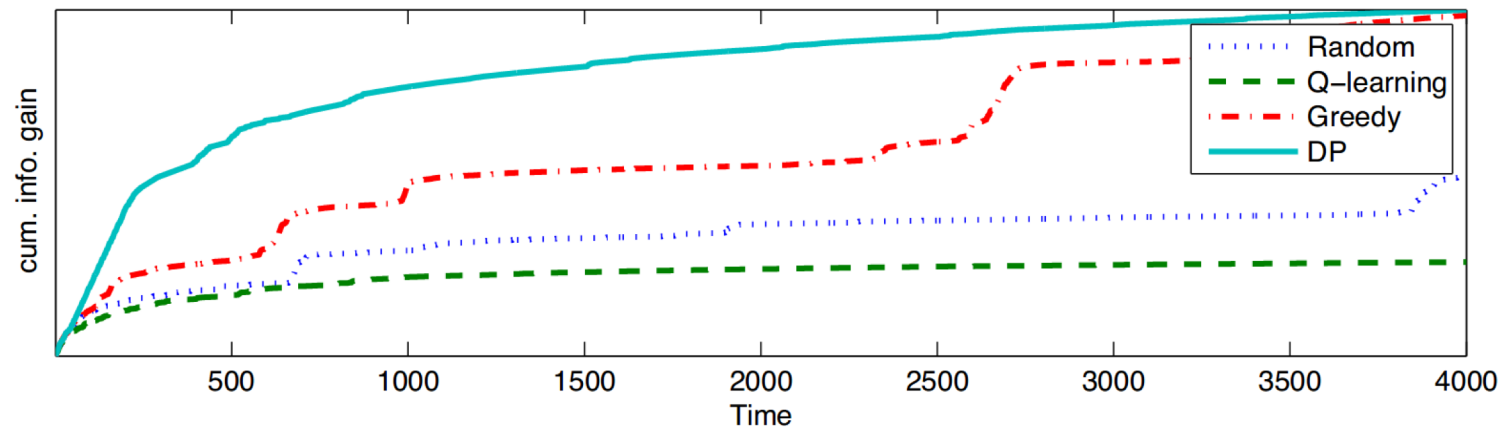
Greedy w.r.t expected one-step info gain

Q-learning using one-step info gain

Policy iteration (Dynamic programming approximation to optimal bayesian exploration)



[Plan] Results



Outline

- Motivation, Related Works and Demo
- Planning to Be Surprised
- Variational Information Maximizing Exploration
- Unifying Count-Based Exploration and Intrinsic Motivation
- Comparisons and Discussion

[VIME] contribution

Dynamic
s model

$$p(s_{t+1} | s_t, a_t; \theta)$$

Variational inference
for posterior
distribution of
dynamics model

$$\phi' = \arg \min_{\phi} \left[\overbrace{D_{\text{KL}}[q(\theta; \phi) \| q(\theta; \phi_{t-1})]}^{\ell(q(\theta; \phi), s_t)} - \mathbb{E}_{\theta \sim q(\cdot; \phi)} [\log p(s_t | \xi_t, a_t; \theta)] \right]$$

$\ell_{\text{KL}}(q(\theta; \phi))$

1-step exploration bonus

$$r'(s_t, a_t, \dot{s}_{t+1}) \leftarrow r(s_t, a_t) + \eta D_{\text{KL}}[q(\theta; \phi'_{n+1}) \| q(\theta; \phi_{n+1})]$$

[VIME] Quantify the information gained

$$\xi_t = (s_1, a_1, s_2, a_2, \dots, s_t)$$

Reminder: PLAN cumulative info gain

$$KL(p(\theta|\xi'_t) || p(\theta|\xi_t))$$

$$D_{\text{KL}}[p(\theta|\xi_t, a_t, s_{t+1}) || p(\theta|\xi_t)]$$

[VIME] Variational Bayes

What's hard?
$$p\left(\underbrace{\theta}_y \mid \xi_t, a_t, \underbrace{s_{t+1}}_x\right) = \frac{p\left(\underbrace{\theta}_y \mid \xi_t, a_t\right) p\left(\underbrace{s_{t+1}}_x \mid \xi_t, a_t; \underbrace{\theta}_y\right)}{p\left(\underbrace{s_{t+1}}_x \mid \xi_t, a_t\right)}$$

Computing posterior for highly parameterized models (e.g. neural networks)

Approximate posterior $q(\theta; \phi) \approx p(\theta \mid \xi_t, a_t, s_{t+1})$ by minimizing $D_{KL}[q(\theta; \phi) \parallel p(\theta \mid \xi_t, a_t, s_{t+1})]$

↓

Minimize negative ELBO
$$\phi' = \arg \min_{\phi} \left[\underbrace{D_{KL}[q(\theta; \phi) \parallel q(\theta; \phi_{t-1})]}_{\ell_{KL}(q(\theta; \phi))} - \overbrace{\mathbb{E}_{\theta \sim q(\cdot; \phi)} [\log p(s_t \mid \xi_t, a_t; \theta)]}^{\ell(q(\theta; \phi), s_t)} \right]$$

[VIME] Optimization for variational bayes

How to minimize negative ELBO?

$$\ell(q(\theta; \phi), s_t)$$

Take an efficient **single second-order (Newton) update step** to minimize negative ELBO:

$$\Delta\phi = H^{-1}(\ell)\nabla_{\phi}\ell(q(\theta; \phi), s_t)$$

[VIME] Estimate 1-step expected info gain

What's hard? $\mathbb{E}_{s_{t+1} \sim \mathcal{P}(\cdot | \xi_t, a_t)} [D_{\text{KL}}[p(\theta | \xi_t, a_t, s_{t+1}) \| p(\theta | \xi_t)]]$

Computing the exact one-step expected info-gain. High-dimensional states

→ Monte-carlo estimation.

$$a_t \sim \pi_\alpha(s_t) \quad s_{t+1} \sim \mathcal{P}(\cdot | s_t, a_t)$$

$$r'(s_t, a_t, s_{t+1}) = r(s_t, a_t) + \eta D_{\text{KL}}[p(\theta | \xi_t, a_t, s_{t+1}) \| p(\theta | \xi_t)]$$

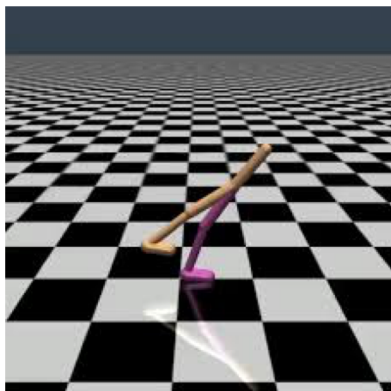
[VIME] Results (Walker-2D)

Dense reward

$$\mathcal{S} \subseteq \mathbb{R}^{20}$$

$$\mathcal{A} \subseteq \mathbb{R}^6$$

RL algorithm: TRPO



Average extrinsic return

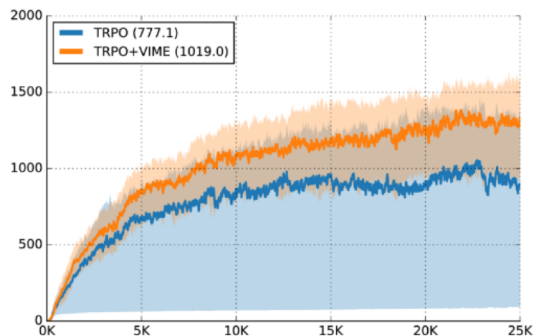


Figure 3: Performance of TRPO with and without VIME on the high-dimensional Walker2D locomotion task.

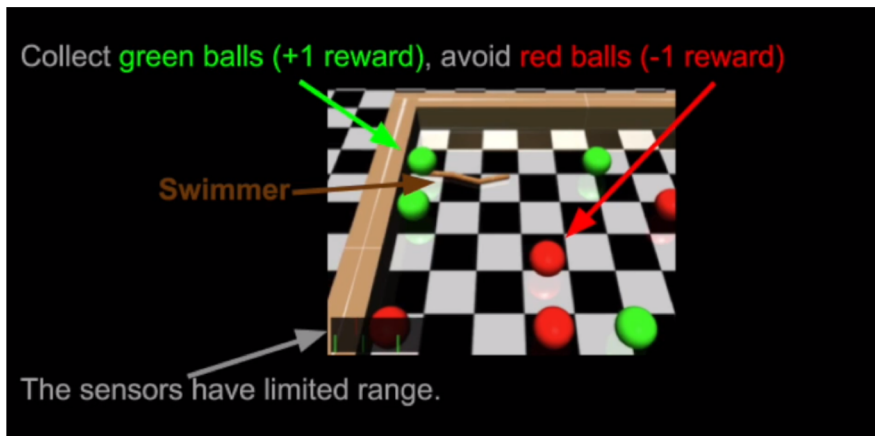
[VIME] Results (Swimmer-Gather)

Sparse reward

$$\mathcal{S} \subseteq \mathbb{R}^{33}$$

$$\mathcal{A} \subseteq \mathbb{R}^2$$

RL algorithm: TRPO



Average extrinsic return

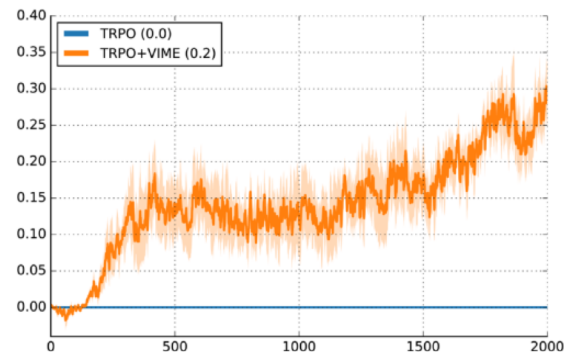
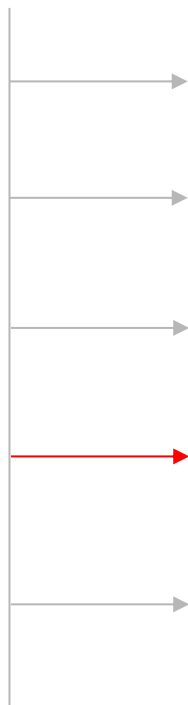


Figure 5: Performance of TRPO with and without VIME on the challenging hierarchical task SwimmerGather.

Outline

- 
- Motivation, Related Works and Demo
 - Planning to Be Surprised
 - Variational Information Maximizing Exploration
 - Unifying Count-Based Exploration and Intrinsic Motivation**
 - Comparisons and Discussion

[CTS] contribution

States Density model

$$\rho_n(x)$$

Pseudo-count

$$\hat{N}_n(x) = \frac{\rho_n(x)(1 - \rho'_n(x))}{\rho'_n(x) - \rho_n(x)}$$

1-step exploration
bonus

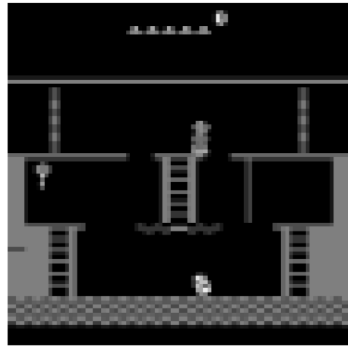
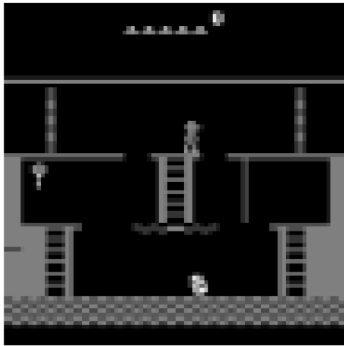
$$R_n^+(x, a) := \beta(\hat{N}_n(x) + 0.01)^{-1/2}$$

[CTS] Count state visitation

Empirical distribution

$$\mu_n(x) := \mu(x; x_{1:n}) := \frac{N_n(x)}{n}$$

← Empirical count



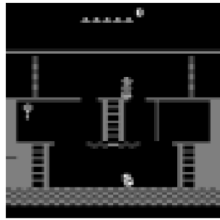
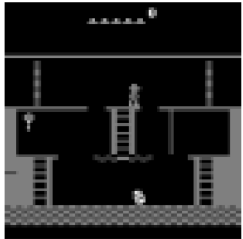
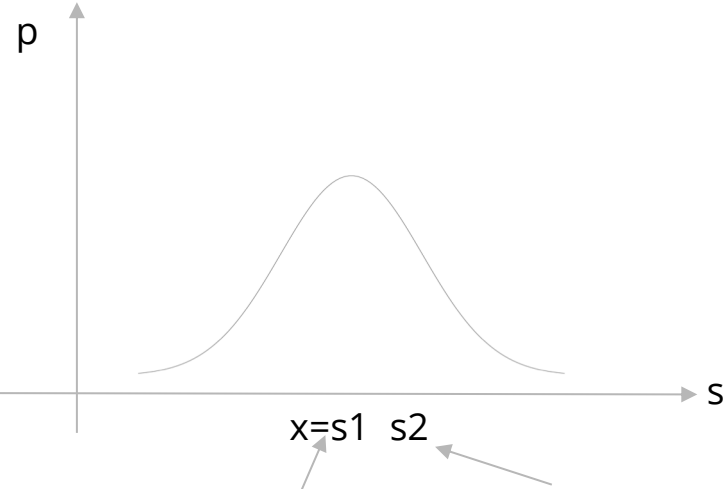
These two are different states!

But we want to increment visitation counts for both when visiting either one.

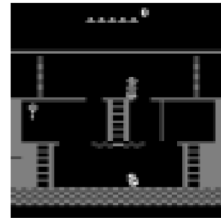
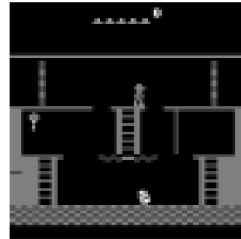
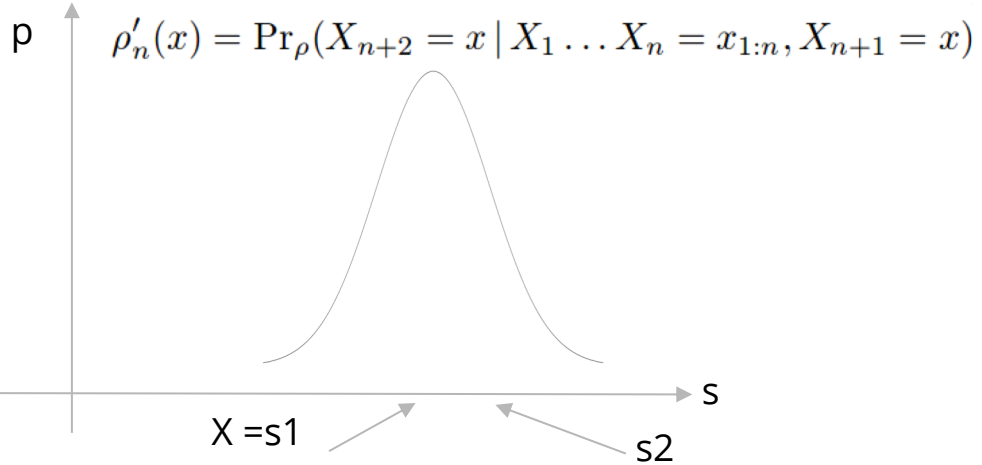
Pixel difference

[CTS] Introduce state density model

$$\rho_n(x) := \rho(x; x_{1:n})$$



$$\rho'_n(x) := \rho(x; x_{1:n}x)$$



How to update CTS density model?

Check the “context tree switching” paper!

<https://arxiv.org/abs/1111.3182>

This was the difficulty of reading this paper as it only shows a bayes rule update for mixture of density models

$$\underbrace{\omega_n(\rho, x)}_{p(y|x)} = \frac{\underbrace{\omega_n(\rho)}_{p(y)} \underbrace{\rho(x; x_{1:n})}_{p(x|y)}}{\underbrace{\int_{\rho \in \mathcal{M}} \omega_n(\rho) \rho(x; x_{1:n}) d\rho}_{p(x)}}$$

Remark: For pixel-cnn density model in “Count-based exploration with **neural density model**”, just **backprop**.

[CTS] Derive pseudo-count from density model

pseudo-count function $\hat{N}_n(x)$
pseudo-count total \hat{n}

Two constraints:
Linear system

$$\rho_n(x) = \frac{\hat{N}_n(x)}{\hat{n}} \quad \rho'_n(x) = \frac{\hat{N}_n(x) + 1}{\hat{n} + 1}$$

Solve linear system

Pseudo-count derived!

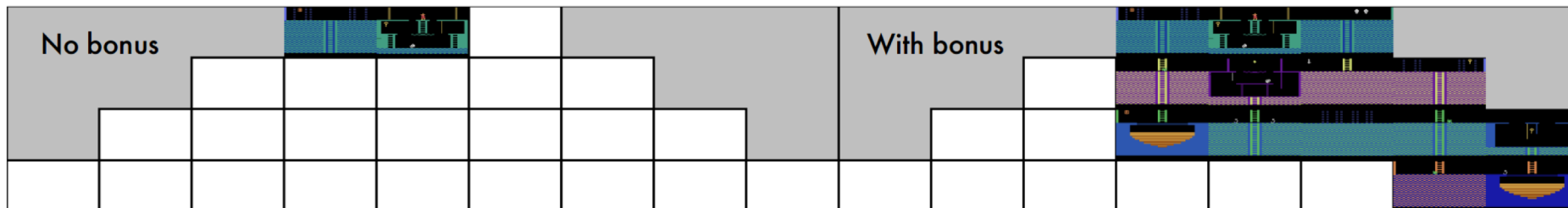
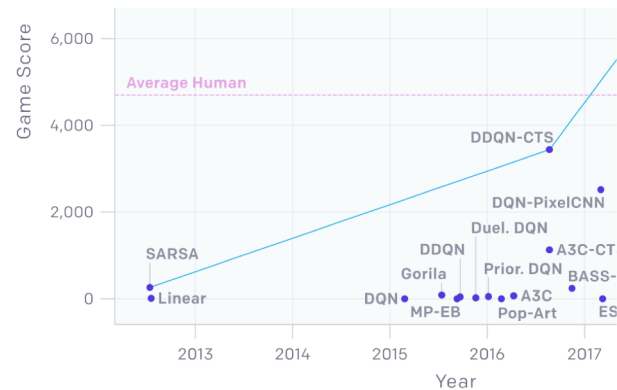
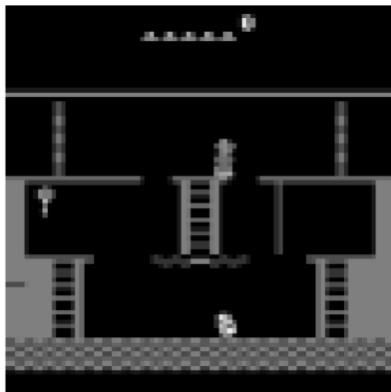
$$\hat{N}_n(x) = \frac{\rho_n(x)(1 - \rho'_n(x))}{\rho'_n(x) - \rho_n(x)}$$

[CTS] Results (Montezuma's Revenge)

State: 84x84x4

Actions: 18

RL algorithm: Double DQN



Outline

- Motivation, Related Works and Demo
- Planning to Be Surprised
- Variational Information Maximizing Exploration
- Unifying Count-Based Exploration and Intrinsic Motivation
- Summary, Comparisons and Discussion

Deriving posterior dynamics model/ density model

PLAN

VIME

CTS

Bayes rule

$$p\left(\underbrace{\theta}_y \mid \xi_t, a_t, \underbrace{s_{t+1}}_x\right) = \frac{p\left(\underbrace{\theta}_y \mid \xi_t, a_t\right) p\left(\underbrace{s_{t+1}}_x \mid \xi_t, a_t; \underbrace{\theta}_y\right)}{p\left(\underbrace{s_{t+1}}_x \mid \xi_t, a_t\right)}$$

Variational inference

$$\phi' = \arg \min_{\phi} \left[\underbrace{D_{\text{KL}}[q(\theta; \phi) \parallel q(\theta; \phi_{t-1})]}_{\ell_{\text{KL}}(q(\theta; \phi))} - \mathbb{E}_{\theta \sim q(\cdot; \phi)} [\log p(s_t \mid \xi_t, a_t; \theta)] \right]$$

Bayes rule

$$\underbrace{\omega_n(\rho, x)}_{p(y|x)} = \frac{\underbrace{\omega_n(\rho)}_{p(y)} \underbrace{\rho(x; x_{1:n})}_{p(x|y)}}{\underbrace{\int_{\rho \in \mathcal{M}} \omega_n(\rho) \rho(x; x_{1:n}) d\rho}_{p(x)}}$$

Derive exploratory policy

Policy trained with the reward augmented by **intrinsic reward**.

[VIME] 1-step Information gain

$$\text{Intrinsic reward} = D_{\text{KL}}[p(\theta|\xi_t, a_t, s_{t+1}) || p(\theta|\xi_t)]$$

[CTS] Pseudo-count

$$\text{Intrinsic reward} = \frac{1}{\sqrt{\hat{N}_n(x)}} \quad \hat{N}_n(x) = \frac{\rho_n(x)(1 - \rho'_n(x))}{\rho'_n(x) - \rho_n(x)}$$

[PLAN] Directly argmax(curiosity Q)

Pseudo-count VS Intrinsic Motivation

Mixture model

$$\underbrace{\omega_n(\rho, x)}_{p(y|x)} = \frac{\underbrace{\omega_n(\rho)}_{p(y)} \underbrace{\rho(x; x_{1:n})}_{p(x|y)}}{\underbrace{\int_{\rho \in \mathcal{M}} \omega_n(\rho) \rho(x; x_{1:n}) d\rho}_{p(x)}}$$

$$\mathbf{IG}_n(x) := \mathbf{IG}(x; x_{1:n}) := \mathbf{KL}(w_n(\cdot, x) \parallel w_n)$$

“Unifying count-based exploration and intrinsic motivations”!

$$\mathbf{IG}_n(x) \leq \hat{N}_n(x)^{-1/2}$$

Limitations & Future Directions

PLAN → Intractable posterior & use dynamics model for expectation
Difficult to be scaled outside Tabular RL.

VIME → Currently maximize sum of 1-step info gain.

CTS → which density model leads to better generalization over states?



Learning rates of policy network VS Updating dynamic model/density model.

Thank you!

(Appendix)

Our derivation for “Additive in expectation”

$$\begin{aligned}
 \mathbb{E}_{p(h''|h')} g(h''|h) &= \mathbb{E}_{p(h''|h')} [KL(p(\theta|h'') || p(\theta|h))] \\
 &= \int \int p(h''|h') p(\theta|h'') \log \frac{p(\theta|h'')}{p(\theta|h)} d\theta dh'' \\
 &= \int \int p(h''|h') p(\theta|h'') \log \frac{p(\theta|h'') p(\theta|h')}{p(\theta|h) \underbrace{p(\theta|h')}} d\theta dh'' \\
 &\quad \text{Mult and div} \\
 &= \int \int p(h''|h') \underbrace{p(\theta|h'')}_{=p(\theta|h'',h')} \log \frac{p(\theta|h')}{p(\theta|h)} d\theta dh'' + \int \int p(h''|h') p(\theta|h'') \log \frac{p(\theta|h'')}{p(\theta|h')} d\theta dh'' \\
 &\quad \text{h'' contains h'} \longrightarrow \\
 &= \int \int p(\theta, h''|h') \log \frac{p(\theta|h')}{p(\theta|h)} d\theta dh'' + \mathbb{E}_{p(h''|h')} g(h''|h') \\
 &= \int \left(\int p(h''|\theta, h') p(\theta|h') dh'' \right) \log \frac{p(\theta|h')}{p(\theta|h)} d\theta + \mathbb{E}_{p(h''|h')} g(h''|h') \\
 &= \int \left(\int p(h''|\theta, h') dh'' \right) p(\theta|h') \log \frac{p(\theta|h')}{p(\theta|h)} d\theta + \mathbb{E}_{p(h''|h')} g(h''|h') \\
 &\quad \text{1} \\
 &= KL(p(\theta|h') || p(\theta|h)) + \mathbb{E}_{p(h''|h')} g(h''|h') \\
 &= g(h'|h) + \mathbb{E}_{p(h''|h')} g(h''|h')
 \end{aligned}$$

